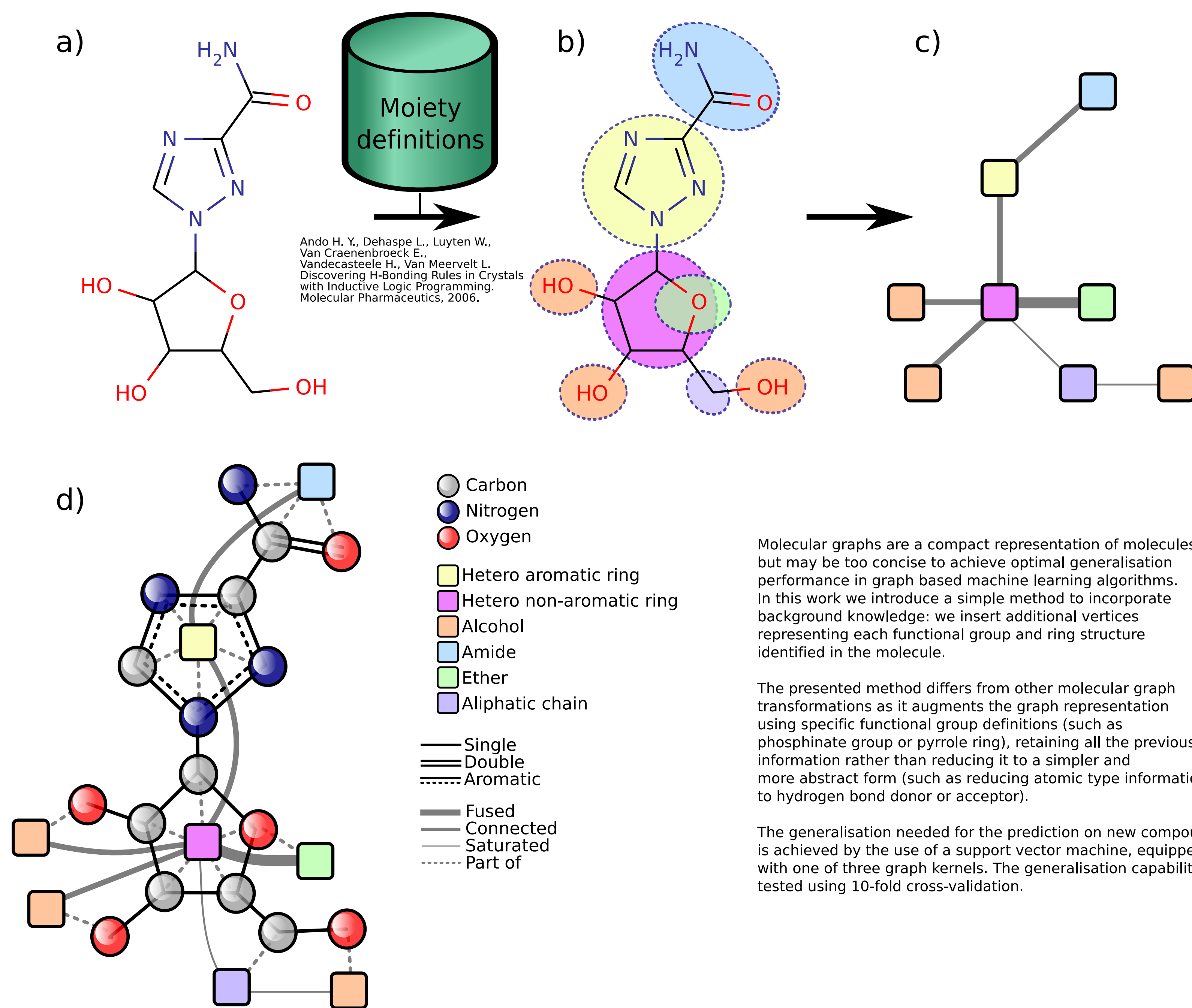


Augmented molecular graph kernel QSARs

Molecular graph augmentation



Existing graph kernels

PAIRWISE DISTANCE KERNEL (PDK)

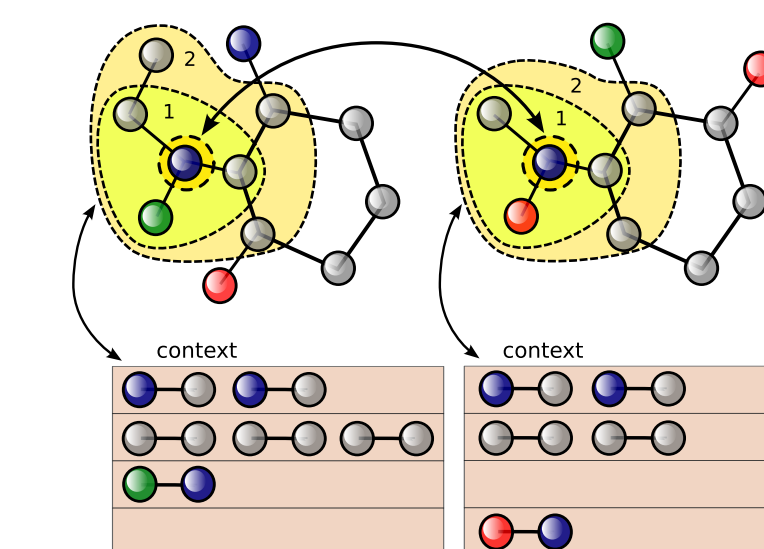
Borgwardt K. M. and Kriegel H.-P. Shortest-Path Kernels on Graphs. In Proc. ICDM, pp. 74-81, 2005.

The idea is to compute the similarity between two graphs by comparing all the respective pairs of vertices annotated with their topological distances. This is achieved by 1) first calculating the shortest path distance between all pairs of vertices using Floyd-Warshall's algorithm, and 2) then computing an all-pairs-shortest-paths kernel on edge walks of length 1 on an appropriately modified graph.

WEIGHTED DECOMPOSITION KERNEL (WDK)

Menchetti S., Costa F., and Frasconi P. Weighted decomposition kernels. In Proc. ICML, pp. 585-592, 2005.

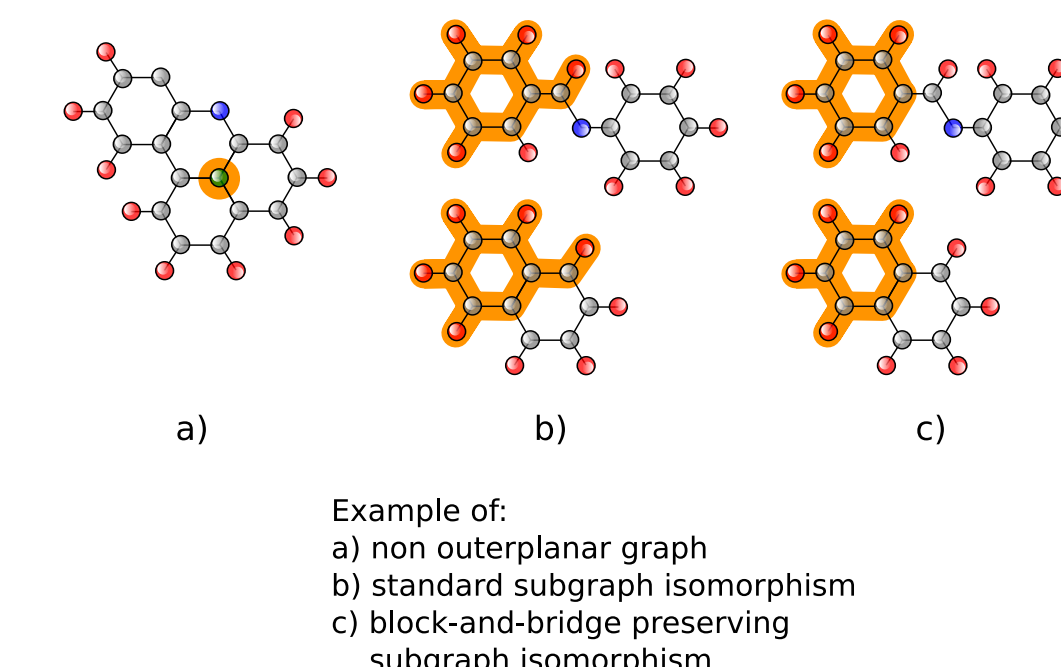
In the WDK, the neighborhood of a given radius is first associated to each vertex in a graph. The WDK is then computed as the product of an exact matching kernel over the vertex label with a kernel over the neighborhood edge multiset. The edge label is augmented with the labels of the endpoints. Among the differences between WDK and NSPDK there are: the single vs. pairwise subgraph approach, and the "soft" similarity match vs. the "hard" isomorphism match of neighbourhood subgraphs.



PAIRWISE MAXIMUM COMMON SUBGRAPHS KERNEL (PMCSK)

Schietgat L., Costa F., Ramon J., and De Raedt L. Maximum common subgraph mining: A fast and effective approach towards feature generation. In Proc. MLG, pp. 1-3, July 2009.

The PMCSK feature space is obtained considering the maximum common subgraph (MCS) between all pairs of instances in the training set. The authors show that, although the computation of the maximum common subgraph in the general case is an NP-hard problem, one can employ a polynomial-time algorithm if only outerplanar graphs are considered in combination with a special case of subgraph isomorphism called block-and-bridge-preserving (BBP) subgraph isomorphism. In addition to the pairwise vs. single subgraph approach, PMCSK differs from NSPDK in the specific type of subgraphs considered (MCSs vs. neighbourhood graphs).



Neighbourhood Subgraph Pairwise Distance Kernel

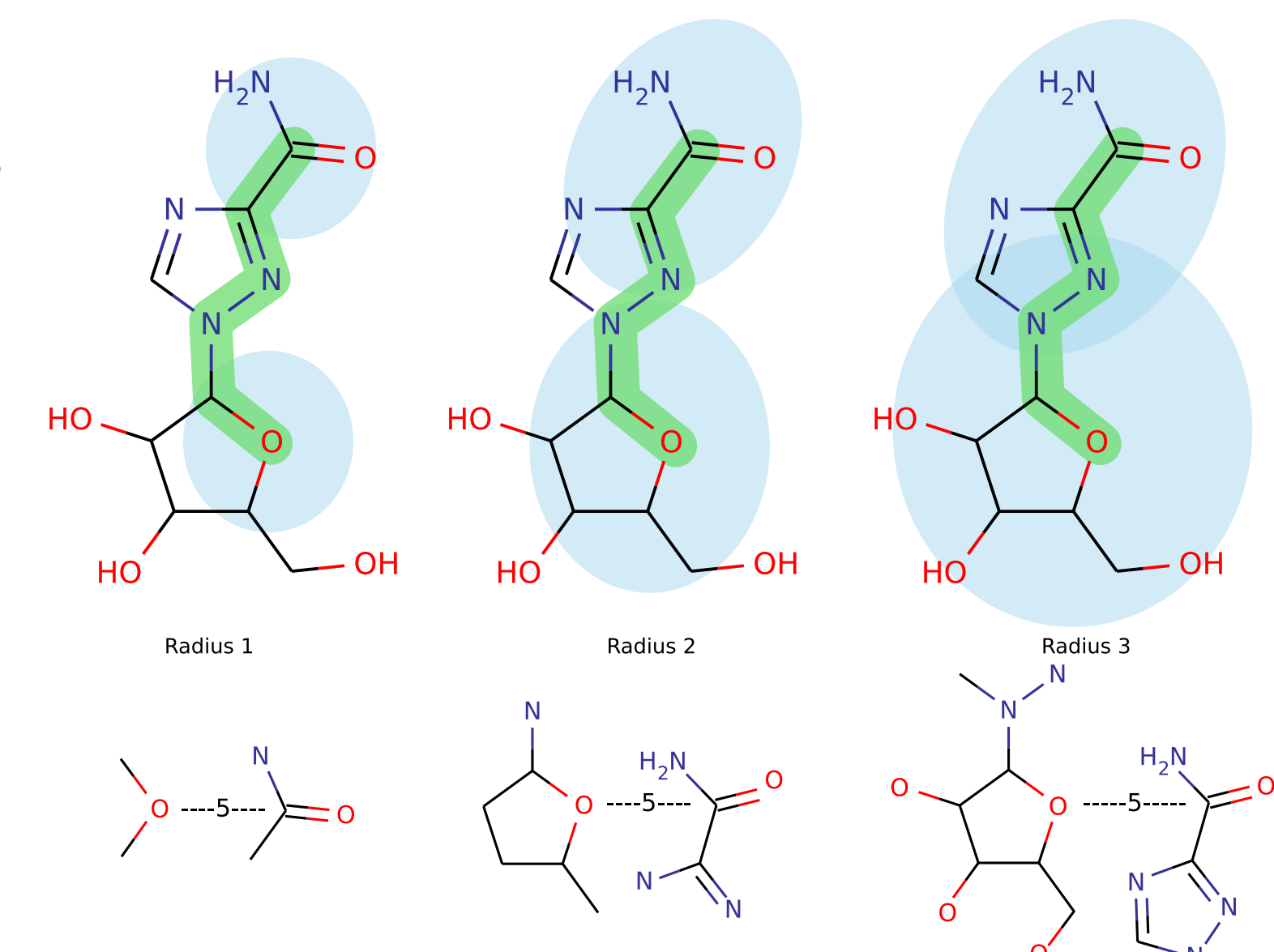
Since the introduction of convolution kernels in (Hausler, 1999), the decomposition approach has been the guiding principle in kernel design for discrete structured objects. The similarity function is obtained by decomposing each object into (possibly overlapping) parts and by devising a local kernel between the subparts.

For over ten years, machine learning researchers have exploited the remarkable property that it is sometimes possible to efficiently compute this type of kernel, even when objects admit an exponential number of decompositions (e.g. through dynamic programming).

However, as the dimension of the feature space associated with the kernel becomes larger, there is an increasing probability that a significant fraction of the feature space dimensions will be poorly correlated with the target function. As a consequence, even when using large margin classifiers, one can fail to obtain models with good generalization performance (Ben-David et al., 2002).

Possible remedies include down-weighting the contribution of larger fragments, or bounding a priori their size. Alternatively, one can try to identify a strong bias relevant to the task at hand, and consider only a selected subset of structures.

We limit the extracted substructures to **pairs of neighbourhood subgraphs, each rooted in one vertex of the given graph and hence efficiently enumerable.**



$$\kappa_{r,d}(G, G') = \sum_{\substack{A_v, B_u \in R_{r,d}^{-1}(G) \\ A_{v'}, B_{u'} \in R_{r,d}^{-1}(G')}} \delta(A_v, A_{v'}) \delta(B_u, B_{u'})$$

$$K(G, G') = \sum_r \sum_d \kappa_{r,d}(G, G').$$

Costa F., De Grave K. Fast neighborhood subgraph pairwise distance kernel. In Proc. International Conference on Machine Learning (ICML) 2010

How fast?

There is no known algorithm of polynomial complexity for exactly matching graphs (isomorphism). We design a fast approximation with very few collisions.

Encoding neighbourhood balls as graph invariant strings allows hashing

ROOTED GRAPH INVARIANT $\mathcal{L}^*(S)$

- Compute all pairwise distances $d(u, v)$
- $\forall u \in V_S$:
 $\mathcal{L}^0(u) \leftarrow \text{concat}(\text{sort}(\{d(u, v), \mathcal{L}(v)\})) \cdot d(u, \text{root})$
- $\forall (u, v) \in E_S$:
 $\mathcal{L}^e(u, v) \leftarrow \mathcal{L}^0(u) \cdot \mathcal{L}^0(v) \cdot \mathcal{L}(u, v)$
- $\mathcal{L}^*(S) \leftarrow \text{concat}(\text{sort}(\mathcal{L}^e))$

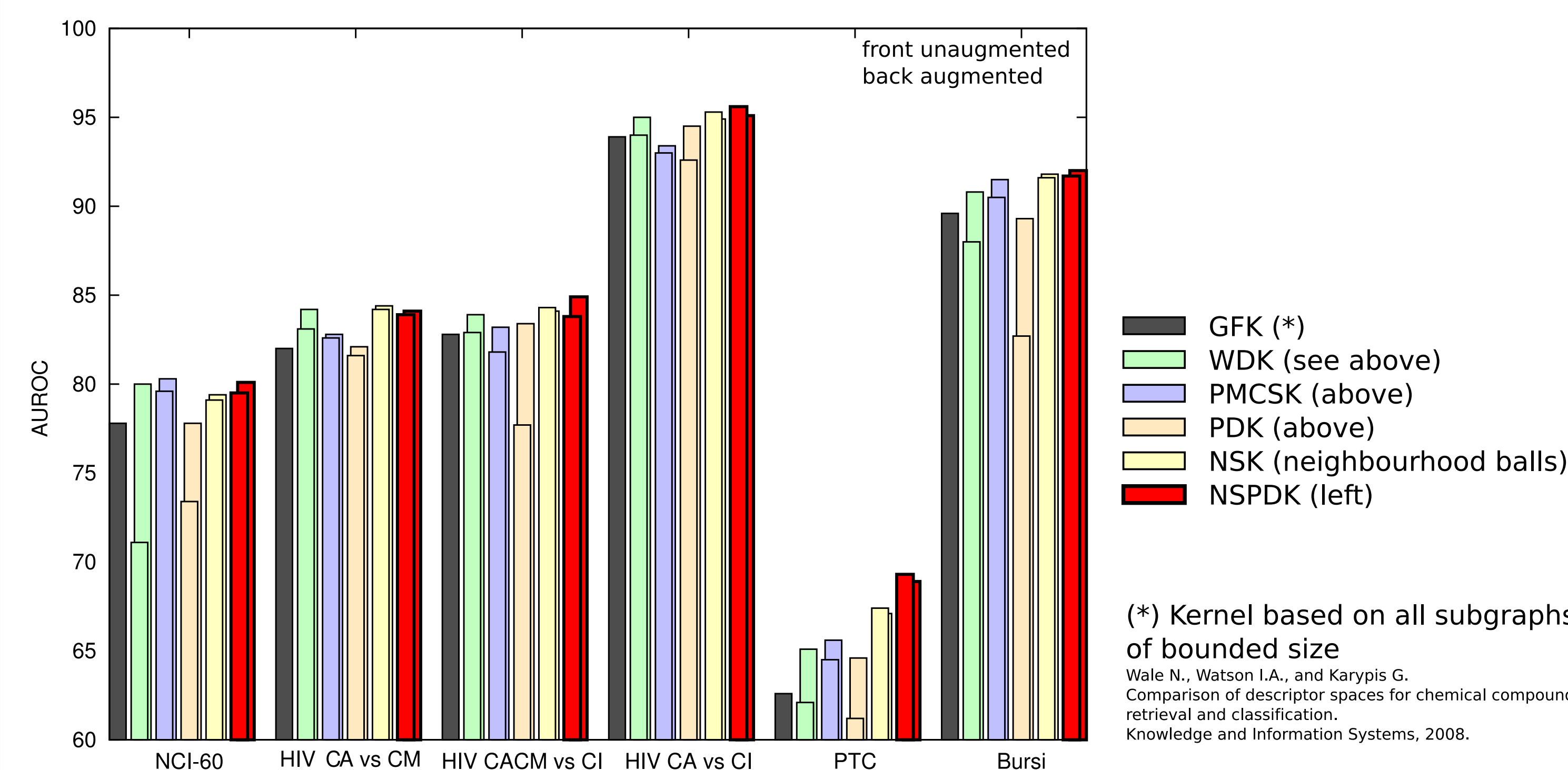
After cryptographic hashing of the features, the kernel can be computed as a sparse vector product.

$$\sim 2 \cdot 10^5 \frac{\text{comparisons}}{\text{core-second}}$$

	NCI-60	HIV	PTC	Bursi
# of mol.	3910	42687	417	4337
Aug. time	$3.5 \cdot 10^2$	$3.4 \cdot 10^3$	$3.4 \cdot 10^3$	$1.2 \cdot 10^2$
GFK(G)	$3.5 \cdot 10^3$	$1.4 \cdot 10^4$	$3.1 \cdot 10^0$	$7.3 \cdot 10^1$
WDK(G)	$1.8 \cdot 10^3$	$1.6 \cdot 10^3$	$8.0 \cdot 10^0$	$1.1 \cdot 10^3$
WDK(G _a)	$2.3 \cdot 10^3$	$2.3 \cdot 10^3$	$1.4 \cdot 10^3$	$1.5 \cdot 10^3$
PMCSK(G)	$2.8 \cdot 10^3$	$3.3 \cdot 10^{4*}$	$6.2 \cdot 10^2$	$3.5 \cdot 10^0$
PMCSK(G _a)	$2.8 \cdot 10^3$	$3.3 \cdot 10^{4*}$	$6.3 \cdot 10^2$	$3.5 \cdot 10^0$
PDK(G)	$4.2 \cdot 10^3$	$3.9 \cdot 10^3$	$1.0 \cdot 10^0$	$3.6 \cdot 10^1$
PDK(G _a)	$7.7 \cdot 10^3$	$4.2 \cdot 10^3$	$2.0 \cdot 10^0$	$5.7 \cdot 10^1$
NSK(G)	$6.2 \cdot 10^3$	$3.1 \cdot 10^3$	$2.8 \cdot 10^0$	$5.1 \cdot 10^1$
NSK(G _a)	$3.5 \cdot 10^2$	$6.0 \cdot 10^3$	$1.4 \cdot 10^3$	$2.0 \cdot 10^2$
NSPDK(G)	$1.2 \cdot 10^2$	$1.0 \cdot 10^2$	$3.4 \cdot 10^0$	$1.1 \cdot 10^2$
NSPDK(G _a)	$4.6 \cdot 10^2$	$1.9 \cdot 10^4$	$1.6 \cdot 10^3$	$2.9 \cdot 10^2$

* MCSs derived only from the 1503 CA-CM molecules.

How accurate?



Conclusions

1. Graph augmentation by and large increases classification performance
2. Increasing the expressiveness of the base kernel, we observe a diminishing return
3. We introduce a fast molecule kernel with excellent generalisation performance